

TITLE OF THE INVENTION

MULTISITE COORDINATION IN SHARED MULTICAST TREES

CROSS-REFERENCE TO RELATED APPLICATIONS

5 This application claims priority from U.S. provisional application serial number 60/214,327 filed on June 27, 2000, incorporated herein by reference.

STATEMENT REGARDING FEDERALLY SPONSORED RESEARCH OR DEVELOPMENT

10 This invention was made with Government support under Grant Nos. F30602-27-0338 and F30602-97-2-0338, awarded by the Air Force Office of Scientific Research, Defense Advanced Research Projects Agency (DARPA). The Government has certain rights in this invention.

REFERENCE TO A MICROFICHE APPENDIX

15 Not Applicable

BACKGROUND OF THE INVENTION

1. Field of the Invention

20 This invention pertains generally to multicast transmission, and more particularly to a protocol to coordinate multipoint groupwork in the IP-multicast framework.

2. Description of the Background Art

Internet computing is gradually migrating from the standard unicast transmission model to multicasting. In the IP multicasting model, a source needs to send a packet only once to the network interface, and multicast routers replicate the packet on its transmission path to multiple receivers. The Internet Group Management Protocol (IGMP) allows a host to join a multicast group by informing its local router to forward multicast traffic for this group to the leaf subnetwork where the host resides. Protocols such as DVMRF, MOSFF, and PIM perform the construction of multicast delivery trees and enable packet forwarding between routers.

With IP multicast, no guarantees are given for reliable or order-preserving delivery of packets, and a message is delivered on a best effort basis to all members of a multicast group. These shortcomings have spurred much research on reliable multicast between end hosts, and on mechanisms to refine IP multicast, such as using addressing information to enable subcasting or anycasting. Subcasting delivers or retrieves data between a source and select members of a multicast group, and anycasting transfers data to anyone member of a group, for example the nearest proxy from a group of servers. While IGMP targets group membership, and multicasting routing protocols are concerned with delivery, no protocols exist to tackle an emerging problem of multisite communication, which is group coordination. This problem surfaces especially for tightly coupled sessions featuring explicit conference membership control.

Group coordination denotes services to support distributed hosts in coordinating their joint activities, including synchronization of flows from different sources, ordered delivery of distributed event information, and the concurrent use of and access to shared resources, referred to as floor control.

5 Early paradigms of group coordination, mutual exclusion, or concurrency control, have been restricted to discrete data domains rather than multimedia contents, using locking to manifest control, and have not been deployed on an Internet scope. Online group coordination in relation to face-to-face meetings has been studied, and those studies show that greediness for time and bandwidth on behalf of non-cooperative users, as well as lack of social cues such as eye contact contribute to coordination problems. While geared toward supporting humans in their interactions, the concept of group coordination also applies to agent-based interaction. For example, methods to mediate resource contention at the user interface level are known. However, many existing systems for online collaboration are proprietary, sparsely documented, and limited to local area networks, or sessions with few users. These include token-based control mechanisms for a shared workspace, as well as systems in which distributed, task-activated floor control serves as a high-level analogy to collision-sensing in channel access.

20 Various other techniques have been developed as well, such as packet-switched voice conferencing in a broadcast and unicast setting, ITU standards 120 and 320 for video conferencing that are circuit-switched and designed for conferences with few users, the MCP coordination protocol between concurrent media flows based on token

passing, frameworks for hierarchical collaboration looking at bandwidth and delay issues, mutual exclusion algorithms operating on a static logical propagation tree in which nodes maintain a pointer to the neighbor node that leads to the current token holder with access to the critical section, adding a dynamic link between a requester and the token holder to avoid backtracked routing of reply messages between a source and target node.

Notwithstanding the foregoing, there is a need to address how routing and end-to-end geometry used for multicasting data can also be used effectively for coordinating the activities among individuals and in sessions of Internet scope. The present invention satisfies that need, as well as others, and overcomes deficiencies in existing group coordination protocols.

BRIEF SUMMARY OF THE INVENTION

In general terms, the present invention is a group coordination method which will be referred to herein as Aggregated Coordination Protocol (ACP). By way of example, and not of limitation, ACP operates on a shared multicast tree, benefiting from the underlying tree structure to store and forward coordination primitives between hosts in different multicast groups on the tree. ACP coordinates distributed activities via message passing, and manifests control by ephemeral permissions rather than actual locks, allowing control over continuous media flows as well as discrete data.

An object of the invention is to provide a coordination mechanism based on IP multicast.

Another object of the invention is to provide a group coordination protocol that can be used for collaborative Internet or Intranet applications.

Another object of the invention is to provide a protocol that is based on a logical control tree, and scales to large groups.

5 Another object of the invention is to provide a protocol that is a standalone mechanism but which can rely on an underlying end-to-end reliable multicast tree.

Another object of the invention is to eliminate the need to build a separate control structure for tracking, routing, withholding, or forwarding control directives.

Further objects and advantages of the invention will be brought out in the following portions of the specification, wherein the detailed description is for the purpose of fully disclosing preferred embodiments of the invention without placing limitations thereon.

BRIEF DESCRIPTION OF THE DRAWINGS

The invention will be more fully understood by reference to the following drawings which are for illustrative purposes only:

FIG. 1 is a diagram showing the structure of a packet header fields for coordination primitive packet according to the present invention.

FIG. 2 is a diagram depicting an example of aggregated coordination protocol operation according to the present invention.

20 FIG. 3 is a chart comparing the average costs to coordinate hosts using conventional unicast, conventional multicast, and aggregated multicast according to the present invention.

DETAILED DESCRIPTION OF THE INVENTION

For illustrative purposes the present invention will now be described with reference to FIG. 1 through FIG. 3.

The present invention, which we refer to as Aggregated Coordination Protocol (ACP), is designed to coordinate multipoint groupwork in the IP-multicast framework by relying on transmission of coordination (control) primitives (directives) between group members across a shared end-to-end multicast tree. ACP benefits from the underlying tree structure to store and forward coordination primitives between hosts in different multicast groups on the tree, coordinates distributed activities via message passing, and manifests control by ephemeral permissions rather than actual locks, thus allowing control over continuous media flows as well as discrete data.

1. Model and Definitions

We define a coordination session $C_S = (H, L)$ as a computer network with hosts H and links $L \subset H \times H$. Communication is by message passing only, and we assume that messages eventually arrive correctly. Each host in a session is client and server for coordination primitives (CP) to other hosts. CPs are issued between hosts to synchronize their joint tasks, to implement causal or total ordering in distributed events, and to mitigate access to shared, but exclusive resources. Coordination management must be aligned with membership operations such as joining, leaving or splitting groups.

The entities involved are users (processes), resources at the various hosts, and the CPs coordinating them. Users assume social roles (e.g., moderator, panelist, student), and both users and agents assume control roles (e.g., controlling who may

work with a resource, or holding the right to work with it). We call the host controlling access and operations of a resource the "floor controller" (FC) for that resource. The host being permitted to access the resource at a given moment is called the "floor holder" (FH). Resources can be located at a specific end host (e.g., camera), in replicated form at every host of a multicast group handling the same resource (e.g., telepointer), or in the network (e.g., voice channel). We distinguish between generic CPs ("grant", "release", "open") and resource or media-specific CPs ("rotate left", "zoom in"). The temporary privilege to work with a multimedia resource is also called "floor". CPs contain the sender id, the receiver ids, the time of creation, the allowed duration, and an optional priority value.

Control over shared resources can be centralized, distributed or a hybrid of both. It can be performed successively by individual session members, partitioned, where various session hosts contribute different control functions, or democratic, where consensus is achieved by negotiation, yielding a new consistent control state. A host holds a floor in his turn for a time interval T , which can be preset or timed out. In our model, we assume that the FC role is associated with a specific host, but it may infrequently rove among hosts. The FH changes at every turn. We assume that each host in a session runs the same coordination protocol, serving requests from other hosts and transmitting requests for resources placed by users or their processes.

2. Aggregated Coordination Protocol (ACP)

2.1 Description

ACP operates on a control tree, comprising three main types of nodes:

(a) holder nodes which host the FH, operating on a resource and being

5 transmission sources;

(b) control nodes which host the FC for a specific resource, and are addressed by other nodes asking for a floor; and

(c) target nodes which receive updates of the operations by a FH.

In addition, nodes on the path from a holder to its targets are referred to as "hop nodes".

Furthermore, "leaf nodes" terminate downward forwarding of control information in the

tree. Coordination Primitives (CPs) are disseminated across a single shared tree

connecting all hosts in a session. The tree corresponds to a single shared

acknowledgment tree for concurrent multicasting, allowing multiple sources instead of

building separate dissemination trees per source and multicast group. The tree can be

a working copy of the reliable multicast tree prepared by a protocol such as Lorax. We

outline tree maintenance for the case that group coordination is deployed in a network

complementing an existing underlying reliable multicast tree.

Various mechanisms for initiation, joining, and advertising of collaborative

sessions exist. We assume that one host, representing itself or a multicast group,

20 initiates a session and advertises the session description, multicast address and media

in use in a session directory. The directory serves as a rendezvous interface, which

allows other hosts to join via a call-up mechanism. The control tree is grown from the

initiating node as the root, and other hosts joining the session are considered first off tree, unicasting request-to-join messages to the inviting root node based on addressing information provided by the session directory. A TTL (time-to-live) field in the join packet restricts the session scope.

5 A successful adoption of a child node to the control tree is confirmed with a bind message. Each newly joined host, as the root of its subtree, locally advertises invitation-to-join messages. It may also be the case that separate subtrees may fission together to create a joint control tree. Each node in the tree has a maximum degree D , which must be high enough to reflect the session structure, but not exceed the capacity of a host to efficiently serve its children. Furthermore, if the tree serves also as a media mixing hierarchy, the permissible height must satisfy the end-to-end delay constraints.

10 Open sessions with dynamic membership may incur frequent joining and leaving, or accidental withdrawal of hosts from the control tree. When the root leaves the control tree, the eldest child in the subtree is the designated new root. Age may be determined by location, joining time, or address labels. A hop node, which lost contact with its
15 parent for a timeout larger than possible congestion delays, must rejoin the tree as described. CPs from a node identified as lost are held at the processing host for a timeout and discarded if the host does not reappear. A lost and rejoined host must resend its pending CPs.

20 Routing of CPs in the tree is performed as follows: if a target node is in the subtree of a node, the CP is routed downward the subtree branch where the target resides, otherwise it is sent upward to its parent node. This operation could be

performed by using only directionality information on where a target is located, as is the known mutual exclusion scheme. However, in contrast to a logical geometry used to propagate critical section requests, our protocol must disseminate CPs involving many resources, media types, and their coordination properties across an infrastructure
5 prescribed by the actual underlying multicast tree.

10
15

We need therefore a more expressive addressing and control mechanism for multisite coordination among multicast groups, because many nodes may assume FC or FH roles during dynamic collaboration session. To this end, ACP assigns recursively and top-down unique prefix labels to each node joining the control tree, such that a child node label contains as prefix the label of its parent. For example, using a binary alphabet, the root receives label 1, its children are numbered 10 and 11, etc. The idea of using binary labels has first been used for routing of instructions in multiprocessor systems and has recently been applied to multicast routing and reliable multicasting. Labels allow nodes to be addressed individually, although being part of one or more possibly overlapping multicast groups. Also, hosts can be placed in the tree independent of their membership in different multicast groups, in contrast to other approaches, which demand that the tree organization reflects group membership.

20

Each CP contains the source's label, the target label(s), a sequence number, a local timestamp, a session wide unique resource id obtained from the session directory, and a floor id, which denotes the temporary access permission or an activity descriptor for a resource. The packet header fields and structure of a CP packet are shown in FIG. 1. *CPid* identifies the coordination primitive and the type of operation,

characterizing various resource modalities. *TTL* indicates the scope of the CP. *Opt* is reserved for priority codings. The timer and sequence numbers tag the CP uniquely in the session and event space. *Checks* is the checksum field. The *Sourceaddr* and *Targetaddrs* fields contain the labels for the sender of a CP and its target nodes.

5 Although the initial source and root will change over the course of a session, the branching properties of the positional tree will not change, when the tree is virtually rehung with shifting of control roles. CPs such as requests for the floor on a resource are propagated across the tree using aggregation. This mechanism corresponds to the solution for the ACK implosion problem, used in reliable multicast to limit the feedback of receivers to a source on lost or corrupt packet. Aggregation limits the control process to local groups, rather than letting CPs that can be satisfied locally, flow back all the way to the holder or controller node. ACP packets are assumed to be transferred reliably, which is guaranteed by the underlying reliable multicast protocol. If ACP performs independently, it needs to supply its own reliability mechanisms.

15 Hosts need to maintain locally the following state: the resource IDs shared from local to the remote hosts, and the remote resource IDs accessed locally, together with their CPids; a state table indicating which resources are locally available or held remotely, together with the id of the remote FH; and a request queue *ReqQ* which collects successive CPs from different hosts (the queue is limited by the number of
20 nodes in the session). If a hop node receives the same CPs from different nodes, it aggregates them into one CP, and checks if a response to these requests can be satisfied locally by polling its own state and the state of neighboring nodes. Otherwise

the composite CP is self-routed up or down in the tree toward the target node(s). This implies that the number of CPs required to coordinate nodes decreases as the request activity increases, because requests are not sent further, if a hop node is reached that already processed the CP. The target address may also contain the name of a
5 multicast group, which is then resolved into its members locally at the primary retriever node for this group. A joining node retrieves the current control state for resources of concern by polling its parent node.

In addition, each node maintains a FIFO queue of pending CPids, identified by the senders' labels. A hop node receiving a floor compares the label of the elected
10 node with the head of the queue, and self-elects if its own label matches the head, or forwards it in the routing procedure outlined above. A control node receiving a request responds immediately to the request by sending back a grant message to the requester, if its local queue is empty, or it appends it to its queue. When control shifts from a node to another one, the pending request queue is transferred to the new control node and its
15 new address is multicast to all groups sharing the associated resource.

Example

Referring to FIG. 2, consider a scenario where three hosts from three different multicast groups MG1, MG2, MG3 must coordinate access to a shared resource for which they are contending. FIG. 2 depicts a snapshot of the protocol operation in a
20 ternary tree.

Assume that hosts 12, 100 and 11 all request the floor held by FC at location 101. All request packets need to be routed along branches of the tree to 101. The

prefix property of the labels allows self-routing of these packets. For example, assume that all hosts are informed that host 101 is FC. Host 12 compares its label with the target label. Its prefix matches (1), but the second identifier indicates that the FC is on subtree 0. The request packet is hence sent upward to host 1, which compares its label with the target, and, detecting that 101 is one of its children, sends the packet to host 10, whose label matches the prefix of 101.

This node performs the same comparison and the packet ultimately arrives at FC, which finally grants the floor to node 12. The forwarding of control directives is aggregated, i.e., multiple requests for the same information from different nodes in the tree are assembled in a hop node in the tree, that lies on the path to the target; and are forwarded combined. This limits control traffic and unnecessary propagation of requests that for example cannot be satisfied at a FC node. FC is hence liberated from the need to communicate with every host in the session, and deals only with relevant requests reaching it from neighbor nodes by self-routing.

2.2 Correctness

Assuming that CPs are sent reliably, a node can receive the privilege associated with a CP only if sends a request to the node controlling the associated resource. A message enacting the CP (e.g., grant-floor) can only be sent by a controlling node. The addressing labels in the tree are unique and a CP is only issued for a specific requesting node. Provided that there is a single holder node for a CP *a priori*, the protocol will hence continue to assign a CP to exactly one node.

If several nodes send a CP and another node holds the privilege associated with it for an infinite time, a deadlock exists. Reasons may be that the wrong node or no node holds the privilege, it is unreachable, or the propagation of a request fails prematurely. Assuming finite FIFO request queues at each node, a request reaching a control node will eventually be served. For a tree of height $h=2$, we have a star-based scheme with $n=D+1$ nodes. A floor will rove between the $D+1$ nodes. Using induction over the height of the tree, we assume that the liveness argument holds for any height l with $2 \leq l < h$. A tree of height h has one additional level, from where CPs can be sent and must be replied to. Because of the acyclic tree geometry, it is impossible that a CP can be outrun by a moving floor privilege. Assume that a node in the additional level of the tree (the root, or a leaf level of a subtree) sends a CP. The addressing semantics of the protocol ensures that messages are self-routed to this level of the tree and will eventually reach the target node. Finally, although privilege passing may progress during a session, a node may starve by being exempted indefinitely from the allocation process. However, the addressing semantics of ACP and the finite FIFO queue forbid that accepted CPs at a control node are ignored indefinitely. Based OR these qualitative arguments we conjecture that ACP is correct.

2.3 Fairness

Fairness refers to the frequency and duration, by which nodes acquire a privilege on the average for a given period, which is for example the session lifetime, or the lifetime of a shared resource in a session. Network latency, geographic distance and location of nodes, or varying host capabilities can all be factors in causing uneven

dissemination patterns of CPs and unfair allocation of floors. Leaf nodes take more time to propagate their requests across the root to a node on the other side of a tree, than nodes just below the root. Shared trees also do not provide shortest paths between a source and its receiver set, which may cause increased latency in CP transfer. It is hence important to establish service policies, which counteract these factors. One simple solution, a "least-recently-served" policy, can be enacted by letting each node maintain a local record of the most recent CPs and their originating nodes. Those nodes are serviced first, which do not appear on the list, or appear last in time or frequency of service.

2.4 Resilience

Previously, we assumed that transfer of CPs and accounting of control information among nodes is failure free. Even we if assume reliable multicasting to ensure that CPs are eventually transferred, the control apparatus may need additional recovery mechanisms to ensure consistency. This applies to regular node failure, control node failure, link failure, or token loss or duplication. Such exceptions can be preempted by redundant dissemination of status information, or by detection of loss and recovery. Regular node or controller failures are typically detected via timeout and recovered with an election protocol, with neighbor nodes providing state updates. Continuation of a split session is possible if the members in each partition agree to continue, e.g., if a quorum exists.

One method to deal with the case that a CP is lost completely or reaches only a subset of nodes is to multicast a CP probe message from a node i to the session

remainder. The response time t_r to receive a response is bound by the maximum time for the probe to traverse the longest link, t_{max} , plus the time t_{ack} for the receiver nodes to send a positive or negative acknowledgment, $t_r = 2t_{max} + t_{ack}$. If the CP is diagnosed as lost, the controller node for the respective floor must regenerate the token and send an
5 update to the session.

2.5 Performance

Attaching positional labels to nodes in a D -ary tree implies an additional storage cost of $\log_2 D$ bits per level in a positional tree of N receivers and height $\log_D N$, i.e., $\lg N$ bits are needed. Using 32-bit labels for designating sources and targets in message headers up to 2^{32} hosts can hence be accommodated. Prefix comparison is cheaper for nodes close to the root due to shorter labels. Serving a CP costs $C_{CP} = C_{req} + C_{resp} + C_{upd}$, comprising the cost to send a request to a control node, receive a response, and multicast an update on the new state. We compare the delay in a unicast, multicast, and aggregated: multicast-communication model under full load (each node sends a
10 CP), assuming that the host processing cost for request, response and update packets is equal and normalized. The average path length between nodes is assumed to be the same for all models. λ represents the individual processing, packetization and transmission delay for each type of packet.

In unicast, the coordination delay incurs $(N - 1)$ requests, replies from control
15 nodes, and updates, where N is the current session size, i.e., $CD_{uc} = 3(N-1)\lambda$. In multicast, $(N - 1)$ nodes send requests, and the control node multicasts one reply and one update back to the session, i.e., $CD_{mc} = (N + 1)\lambda$. In aggregated multicast, CPs are

handled within multicast groups and only the root of a group forwards a composite request to its parent, or responds to group-local requests, if it holds the information locally. With K groups we have on the average $G = \left\lceil \frac{N}{K} \right\rceil$ members per group, and per group there are G requests inside a group, K aggregated requests sent to a control node from all groups, and one multicast response and update, i.e.,

$$CD_{amc} = ((G - 1) + (K - 1) + 2)\lambda = (G + K)\lambda.$$

FIG. 3 shows the average cost to coordinate hosts in sessions up to size $N = 1000$, clustered into $K = N/10$ groups, and with normalized transmission delay λ . It elicits the benefits of aggregated multicast coordination.

Although the description above contains many specificities, these should not be construed as limiting the scope of the invention but as merely providing illustrations of some of the presently preferred embodiments of this invention. Therefore, it will be appreciated that the scope of the present invention fully encompasses other embodiments which may become obvious to those skilled in the art, and that the scope of the present invention is accordingly to be limited by nothing other than the appended claims, in which reference to an element in the singular is not intended to mean "one and only one" unless explicitly so stated, but rather "one or more." All structural, chemical, and functional equivalents to the elements of the above-described preferred embodiment that are known to those of ordinary skill in the art are expressly incorporated herein by reference and are intended to be encompassed by the present claims. Moreover, it is not necessary for a device or method to address each and every

problem sought to be solved by the present invention, for it to be encompassed by the present claims. Furthermore, no element, component, or method step in the present disclosure is intended to be dedicated to the public regardless of whether the element, component, or method step is explicitly recited in the claims. No claim element herein is
5 to be construed under the provisions of 35 U.S.C. 112, sixth paragraph, unless the element is expressly recited using the phrase "means for."

10
15
20
25
30
35
40
45
50
55
60
65
70
75
80
85
90
95
100
105
110
115
120
125
130
135
140
145
150
155
160
165
170
175
180
185
190
195
200
205
210
215
220
225
230
235
240
245
250
255
260
265
270
275
280
285
290
295
300
305
310
315
320
325
330
335
340
345
350
355
360
365
370
375
380
385
390
395
400
405
410
415
420
425
430
435
440
445
450
455
460
465
470
475
480
485
490
495
500
505
510
515
520
525
530
535
540
545
550
555
560
565
570
575
580
585
590
595
600
605
610
615
620
625
630
635
640
645
650
655
660
665
670
675
680
685
690
695
700
705
710
715
720
725
730
735
740
745
750
755
760
765
770
775
780
785
790
795
800
805
810
815
820
825
830
835
840
845
850
855
860
865
870
875
880
885
890
895
900
905
910
915
920
925
930
935
940
945
950
955
960
965
970
975
980
985
990
995